

MARKOV KATANA

1 **Markov Katana: a Novel Method for Bayesian Resampling of**
2 **Parameter Space Applied to Phylogenetic Trees**

3

4 Stephen T. Pollard¹, Kenji Fukushima¹, Zhengyuan O. Wang², Todd A. Castoe³ and David D.
5 Pollock^{1*}

6 ¹Department of Biochemistry and Molecular Genetics, University of Colorado School of
7 Medicine, Aurora, CO 80045

8 ²Washington University School of Medicine, St. Louis, MO 63110

9 ³University of Texas Arlington, TX 90238

10 *Correspondence to: David.Pollock@ucdenver.edu, 12801 E 17th Ave, MS 8101, Aurora, CO
11 80045; 303-724-3234

12

13

14 ABSTRACT

15 Phylogenetic inference requires a means to search phylogenetic tree space. This is usually
16 achieved using progressive algorithms that propose and test small alterations in the current tree
17 topology and branch lengths. Current programs search tree topology space using branch-
18 swapping algorithms, but proposals do not discriminate well between swaps likely to succeed or
19 fail. When applied to datasets with many taxa, the huge number of possible topologies slows
20 these programs dramatically. To overcome this, we developed a statistical approach for proposal
21 generation in Bayesian analysis, and evaluated its applicability for the problem of searching
22 phylogenetic tree space. The general idea of the approach, which we call ‘Markov katana’, is to
23 make proposals based on a heuristic algorithm using bootstrapped subsets of the data. Such
24 proposals induce an unintended sampling distribution that must be determined and removed to
25 generate posterior estimates, but the cost of this extra step can in principle be small compared to
26 the added value of more efficient parameter exploration in Markov chain Monte Carlo analyses.
27 Our prototype application uses the simple neighbor-joining distance heuristic on data subsets to
28 propose new reasonably likely phylogenetic trees (including topologies and branch lengths). The
29 evolutionary model used to generate distances in our prototype was far simpler than the more
30 complex model used to evaluate the likelihood of phylogenies based on the full dataset. This
31 prototype implementation indicates that the Markov katana approach could be easily
32 incorporated into existing phylogenetic search programs and may prove a useful alternative in
33 conjunction with existing methods. The general features of this statistical approach may also
34 prove useful in disciplines other than phylogenetics. We demonstrate that this method can be
35 used to efficiently estimate a Bayesian posterior.

36 Key words: phylogenetics, tree search, Bayesian, bootstrap

37 INTRODUCTION

38 Phylogenetic inference has long played a pivotal role in molecular evolution and evolutionary
39 genomics (e.g. Felsenstein 2004; Vonk 2013; Fukushima 2017). It provides unique information
40 about gene and protein interactions (Wang 2005; Hackett 2007; Reyes-Prieto 2007; Craig 2007)
41 and is critical for detecting adaptive bursts and functional divergence (e.g. Castoe 2008; Castoe
42 2009). Despite its importance, phylogenetic inference is difficult partly because searching tree
43 space is an NP-hard problem (Bodlaender 1992; Brocchieri 2001). Distance-based methods such
44 as neighbor-joining (NJ; (Saitou 1987)) are fast and often provide good approximate results but
45 are considered less reliable than the computationally expensive (Hershkovitz 1998; Takahashi
46 2000; Whelan 2001) likelihood-based methods (maximum likelihood, ML, and Bayesian or
47 posterior probability, PP). While distance methods generate a single tree using heuristic
48 approaches, likelihood methods must search tree space, generally by running an optimization
49 scheme or Markov chain Monte Carlo (MCMC). Tree space is often searched using various
50 forms of branch swapping (Felsenstein 1981; Huelsenbeck 1997, 2001; Sullivan 2005;
51 Anisimova 2006). A cautious approach to interpreting results from traditional branch-swapping
52 algorithms is warranted, particularly for trees with sequences from many taxa (Mossel 2005).

53 The principle confounding effect in phylogenetic inference is that multiple substitutions
54 may occur at the same site. Distance-based methods are inferior to likelihood-based methods in
55 accurately inferring multiple substitutions (Felsenstein 1984; Huelsenbeck 1996; Xia 2006).
56 Distance-based methods are also far more strongly biased by long-branch attraction and cannot
57 fully incorporate the advantages of site-specific models of evolution (Huelsenbeck 1995, 1997;
58 Pollock 1998). Another major class of phylogenetic analysis, based on the principle of maximum
59 parsimony, will not be considered here because parsimony methods are far slower than distance

MARKOV KATANA

60 methods, and they do not accurately model evolutionary processes despite having the same
61 biases and inaccuracies as distance methods. The computational limitations of likelihood-based
62 methods become far more severe with large amounts of sequence data from highly diverse sets
63 of organisms (Pollock 2000; Sanderson 2003; A. J. de Koning 2010). For example there are
64 2.75×10^{76} possible topologies relating 50 taxa (Felsenstein 2004), making exhaustive approaches
65 impossible. Branch-and-bound searches can reduce the tree space to be examined for smaller
66 trees but are insufficient for large datasets because the number of tree topologies is still too large
67 (Hendy 1982). Thus, heuristic searches must be used for large trees, evaluating trees that are
68 proximal to reasonably likely trees that have already been found. These searches are currently
69 often performed using branch-swapping algorithms such as nearest-neighbor interchange (NNI),
70 subtree pruning and regrafting (SPR) and tree bisection and reconnection (TBR) (e.g. Ronquist
71 2003; Salemi 2009). The number of NNI, SPR, and TBR neighbors of any topology increase
72 respectively as linear, quadratic, and cubic functions of the number of taxa, and the trees
73 proposed are not necessarily of similar likelihood to the known tree. Therefore, many highly
74 improbable trees are evaluated in branch-swapping algorithms, and the correct solution is not
75 guaranteed due to the presence of local optima in tree space (Mossel 2005). Branch length
76 optimization (or posterior equilibration) must also be performed after branch swapping and is an
77 additional source of computational cost.

78 Several heuristic approaches have been developed to release tree searches from local
79 optima. Ratchet methods employ multiple initial trees perturbed by bootstrap resampling to
80 ensure a less-overlapping tree space in subsequent optimizations using branch swapping (Nixon
81 1999; Vos 2003). The partial stepwise addition (PSA) approach enables escape from local
82 optima by removing some taxa during the topology search (Whelan 2007). Simulated annealing

MARKOV KATANA

83 (SA; Kirkpatrick 1983) and Metropolis-coupled Markov chain Monte Carlo (MCMCMC; Geyer
84 1991) manipulate a likely range of proposed tree acceptances in a single heuristic search or in
85 multiple interacting chains, respectively. Genetic algorithms (GAs) simulate the population
86 dynamics of tree topologies using likelihood as a fitness parameter (Matsuda 1995). These
87 methods outperform simple heuristic searches in at least some contexts. All approaches listed
88 above employ branch swapping to explore tree space and therefore suffer from inefficiency due
89 to the decoupling of topology proposals from the likelihoods of the topologies.

90 Here we consider whether the beneficial features of Bayesian analyses under relatively
91 complex models can be profitably combined with the speed of distance methods based on
92 relatively simple models. The key to our approach is that rather than using branch swapping to
93 explore phylogenetic tree space, distance-based trees predicted from partially sampled sequences
94 are used. We use Markov chain Monte Carlo (MCMC) and a Metropolis-Hastings algorithm in
95 which new steps in the chain are proposed based on bootstrap resampling a proportion of the
96 current sequence sample. Heuristic phylogenetic trees based on the new sample are created using
97 NJ and the likelihoods of the new trees are evaluated using the full sequence dataset and the
98 mtMam model (Yang 1998). The unwanted sampling distribution induced by the NJ proposal
99 mechanism is estimated by running the proposal mechanism without calculating the likelihoods
100 of the proposed trees. The posterior is then corrected for this sampling distribution. We evaluated
101 the effect of different site sample sizes used to generate the NJ trees (sample size) and different
102 resample proportions (jump size).

103 MATERIALS AND METHODS

104 *Mitochondrial Sequences*

105 The 495 amino acid COI 1249-taxon mitochondrial gene alignment from Goldstein et al
106 was used (Goldstein 2015). 10 taxa were arbitrarily selected from the alignment to use for testing
107 and are shown in Table 1.

108 *Program Details*

109 A Perl program, *MarkovKatana*, was written to implement the Markov chain
110 bootstrapping algorithm. *MarkovKatana* takes multiple sequence alignments in the fasta format
111 and outputs phylogenetic trees in the Newick format, along with likelihood values. Another
112 program *Forest* was written to analyze the trees generated by *MarkovKatana* to calculate tree
113 and branch frequencies. *MarkovKatana* and *Forest* were tested on and are compatible with
114 current Unix-based operating systems as well as Windows. The program *PAML* was used to
115 calculate the likelihoods for the trees using the entire alignment of 495 amino acids (Yang 2007).

116 *Branch Prior Calculations*

117 Branch priors were calculated as

118
$$P(B_1 | N, s_l) = \frac{T_s^r * T_{N-s}^r}{T_N^u}, \quad (1)$$

119 where T_x^r and T_x^u are respectively the number of possible rooted and unrooted topologies with x
120 taxa, N is the total number of taxa being evaluated, and s is the smaller number of taxa that are
121 segregated on one side or the other of branch B_b (Pickett 2005).

122 *Modifying Implementation of NJ in Markov Katana to Improve Branch*

123 *Length Estimation*

124 In initial runs, the NJ algorithm often generated unrealistically short branches, so to
125 counteract this we lengthened the shortest branches by adding a random number from 0 to 2
126 substitutions (a branch length increase of 0 to 2/495). This limited the effect of these implausibly
127 short branches in the proposal mechanism. Short branches were still possible, but extremely
128 short branches were not as likely to be proposed.

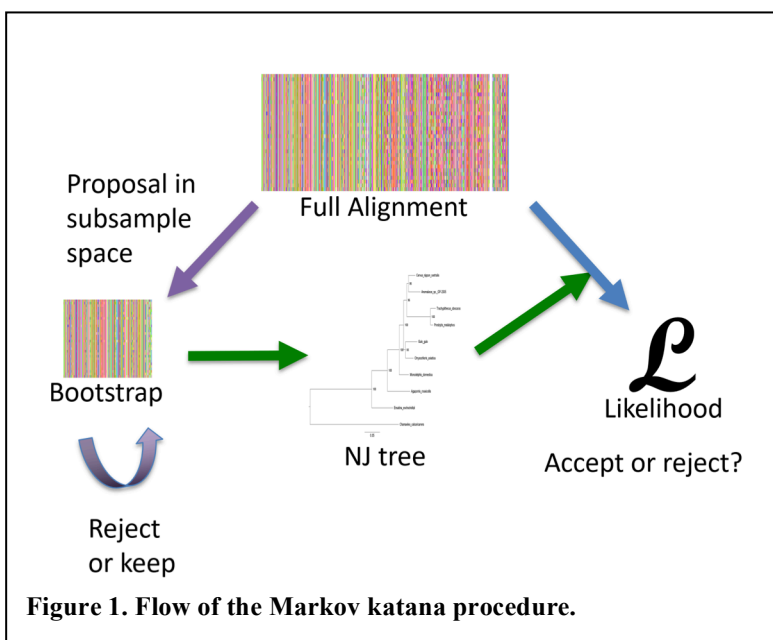
129 RESULTS

130 *Details of the Markov Katana Implementation*

131 A bootstrap sampling procedure (Felsenstein 1985; Zharkikh 1995) was employed to
132 sample sites in the alignment that were then used to calculate distance matrices. Although
133 complete and partial bootstrapping has been used extensively in phylogenetic studies to evaluate
134 branch support and tree confidence (Efron 1996; Alfaro 2003), we used it solely to generate a
135 broad distribution of reasonably likely trees based on the NJ heuristic. Note that while partial
136 sampling is more common when employing the related jackknife approach, bootstrapping
137 approaches such as that employed here sample with replacement, rather than without
138 replacement as in the jackknife. Depending on the number of sites sampled (the sample size),
139 trees produced from partial sequence samples can be quite different from the ML tree of the
140 entire alignment and considerably less likely (Castoe 2009). Evaluating the posterior distribution
141 with an importance sampling approach using these trees is not feasible because the extreme
142 variation in likelihoods among trees means that a few trees would dominate the weighted

MARKOV KATANA

143 importance sampling average (Kuhner 1995). Instead, it is necessary to use a progressive
144 Markov chain approach to evaluate the posterior, such as the Metropolis-Hastings algorithm
145 (Hastings 1970), in which the proposed sample depends on the current sample (Fig. 1). Only a
146 fraction of sites is resampled in each generation of the chain. The NJ tree generated from the
147 proposed sample updates both
148 branch lengths and topology
149 simultaneously, and the
150 likelihood of this proposal was
151 then calculated on the full
152 alignment. The number of sites
153 resampled was uniform
154 randomly chosen up to some
155 maximum, which we will call
156 the ‘jump size’.



157 *Posterior Calculations*

158 To obtain the posterior, the uncorrected distribution of trees after the initial Markov
159 katana (MK) run must be corrected for the bias induced by the proposal mechanism. In these
160 runs, the sample size as a fraction, f , and the jump size, j , were variable parameters and differed
161 among runs as specified. For a given sampled generation, k , the alignment sample at that
162 generation produced a NJ genealogy, G_k , with topology, T_i . The proportion of times that each
163 different topology was produced by the chain out of K sampled generations in the chain is an
164 estimator of the uncorrected posterior for a given sample size, f , or

MARKOV KATANA

165
$$\hat{U}_f(T_i) = \frac{1}{K} \sum_{k=1}^K \delta(T_i, G_k) \quad (2)$$

166 where $\delta(T_i, G_k)$ is a delta function equal to 1 if G_k has topology T_i and otherwise 0. To obtain
 167 the corrected topology posterior, $C(T_i)$, we first estimate the topology sampling bias $\hat{\beta}_f(T_i)$
 168 induced by NJ proposals with sampling fraction f , by sampling K' genealogies from a separate
 169 chain in which all proposals are accepted, to obtain

170
$$\hat{\beta}_f(T_i) = \frac{1}{K'} \sum_{k'=1}^{K'} \delta(T_i, G_{k'}) \quad (3)$$

171 We note that this procedure is identical to obtaining a NJ partial bootstrap, but by running
 172 the Markov chain with a given jump size we can obtain the connectedness among topologies,
 173 providing a natural topological distance measure.

174 We then recognize that

175
$$U_f(T_i) \propto \beta_f(T_i) \int L(G_k) P(G_k) \cdot \delta(T_i, G_k) = \beta_f(T_i) C(T_i) \quad (4)$$

176 where $L(G_k)$, $\beta_f(G_k)$, and $P(G_k)$, are the likelihood, the genealogy sampling bias induced by
 177 NJ proposals with sampling fraction f , and the prior, respectively. Here we assume a flat prior
 178 across all tree topologies. The next step is to divide the uncorrected topology posterior by the
 179 sampling distribution induced by the proposals to obtain

180
$$\hat{C}(T_i) \propto \frac{\hat{U}_f(T_i)}{\hat{\beta}_f(T_i)} \quad (5)$$

181 We normalized the corrected posteriors by dividing by the sum of all corrected posteriors over
 182 all topologies sampled.

183 It may sometimes be useful and possibly more accurate to calculate branch (a.k.a. a

MARKOV KATANA

184 species bi-partition, or edge) posteriors directly over the sample of trees,

$$185 \quad \tilde{U}_f(B_l) = \frac{1}{K} \sum_k \delta(B_l, G_k) \quad (6)$$

186 where $\delta(B_l, G_k)$ is a delta function equal to 1 if G_k has branch B_l , and otherwise 0. The \sim
187 symbol indicates that the branch uncorrected posteriors were calculated directly. In this case, it is
188 necessary to appropriately adjust for the sampling distribution on the branch induced by
189 topological constraints (Pickett 2005), which is contained in both $\tilde{U}_f(B_l)$ and a similarly
190 obtained

$$191 \quad \tilde{\beta}_f(B_l) = \frac{1}{K'} \sum_{k'=1}^{K'} \delta(B_l, G_{k'}) \quad (7)$$

192 This prior is put back into the posterior calculation as

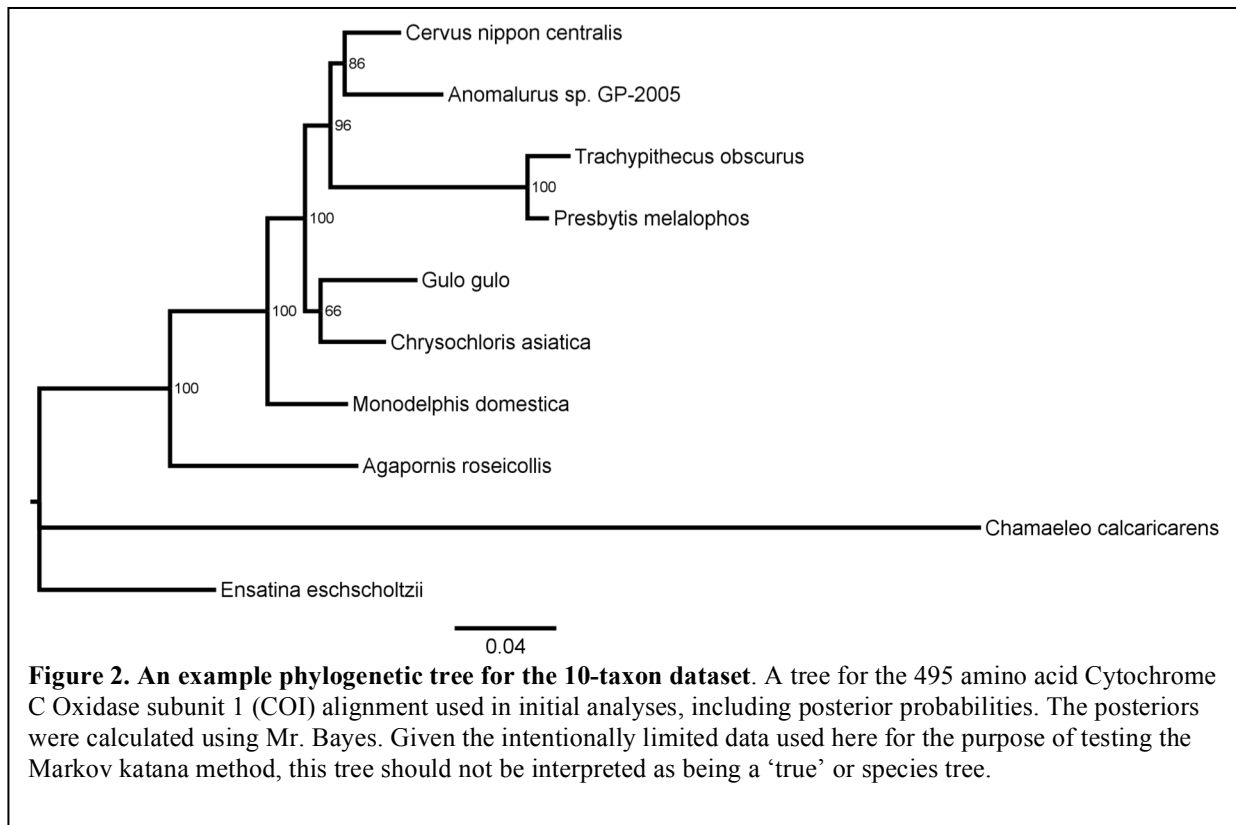
$$193 \quad \hat{C}(B_l) \propto \frac{\tilde{U}_f(B_l)}{\tilde{\beta}_f(B_l)} * P(B_l | N, s_l) \quad (8)$$

194 where $P(B_l | N, s_l)$ is the prior probability of branch B_l induced by topological structures, N is
195 the total number of extant species in the tree, and s_l is the smaller number of species that are
196 partitioned to one side of branch B_l . $P(B_l | N, s_l)$ can be calculated directly (see Methods).

197 *Implementation of the Markov katana*

MARKOV KATANA

198 We began by analyzing a 10-taxon Cytochrome C Oxidase subunit 1 (COI) amino acid
199 alignment (495 residues) that was chosen so that there would be a moderate level of topological
200 uncertainty in the posterior (Fig. 2). Preliminary evaluations indicated that NJ trees on



201 bootstrapped data have a distribution of topologies that are relatively similar among distance
202 types (Sup. Fig. S1). Although there is considerable noise to the estimates for very small
203 frequencies, and there is a slight shift towards higher frequencies with the Markov katana
204 difference NJ, overall the two measures have a nearly linear relationship. This gave us
205 confidence that NJ trees based on differences rather than corrected distances might be
206 sufficiently accurate for our purposes, so to keep the NJ calculations as simple and fast as
207 possible for initial testing, distances were generated using the simple difference matrix. The
208 likelihood of the proposed tree topology and branch lengths were then evaluated using the
209 mtMam substitution rate model (Yang 1998) on the entire sequences. Continuing to keep things

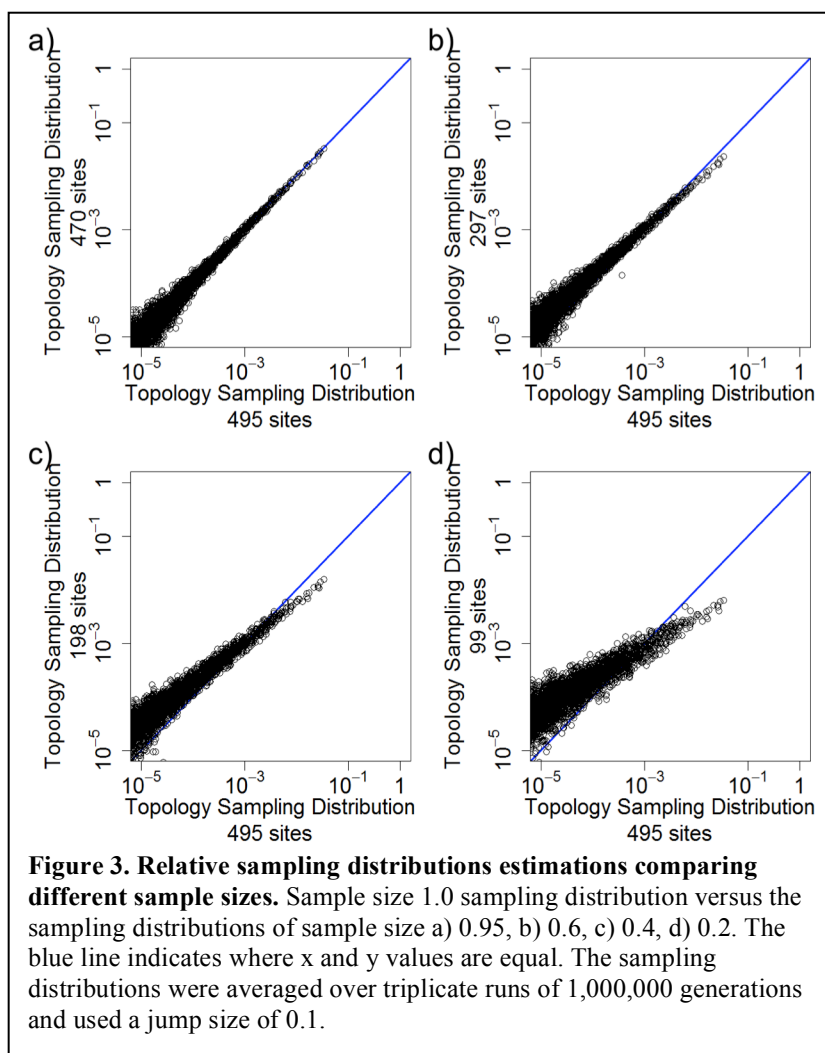
MARKOV KATANA

210 simple for initial testing, we used a flat prior, although we imagine that most future
211 implementations will want to incorporate other priors here, such as the commonly used
212 exponential priors on branch lengths (Yang 2005).

213 To understand the differences in topology sampling bias estimates obtained using
214 different sample sizes, f , Markov katana was run with sample fractions ranging from 100% (495
215 sites) down to 20% (99 sites). The topology sampling biases for smaller f become somewhat
216 more even, with the least frequent topologies about 10x more frequent for $f=20\%$ than for
217 $f=100\%$ (Fig. 3). At the same time, the number of topologies with sampling probabilities greater
218 than 10^{-6} increased from
219 5,975 for $f=100\%$ to
220 21,198 for $f=20\%$.

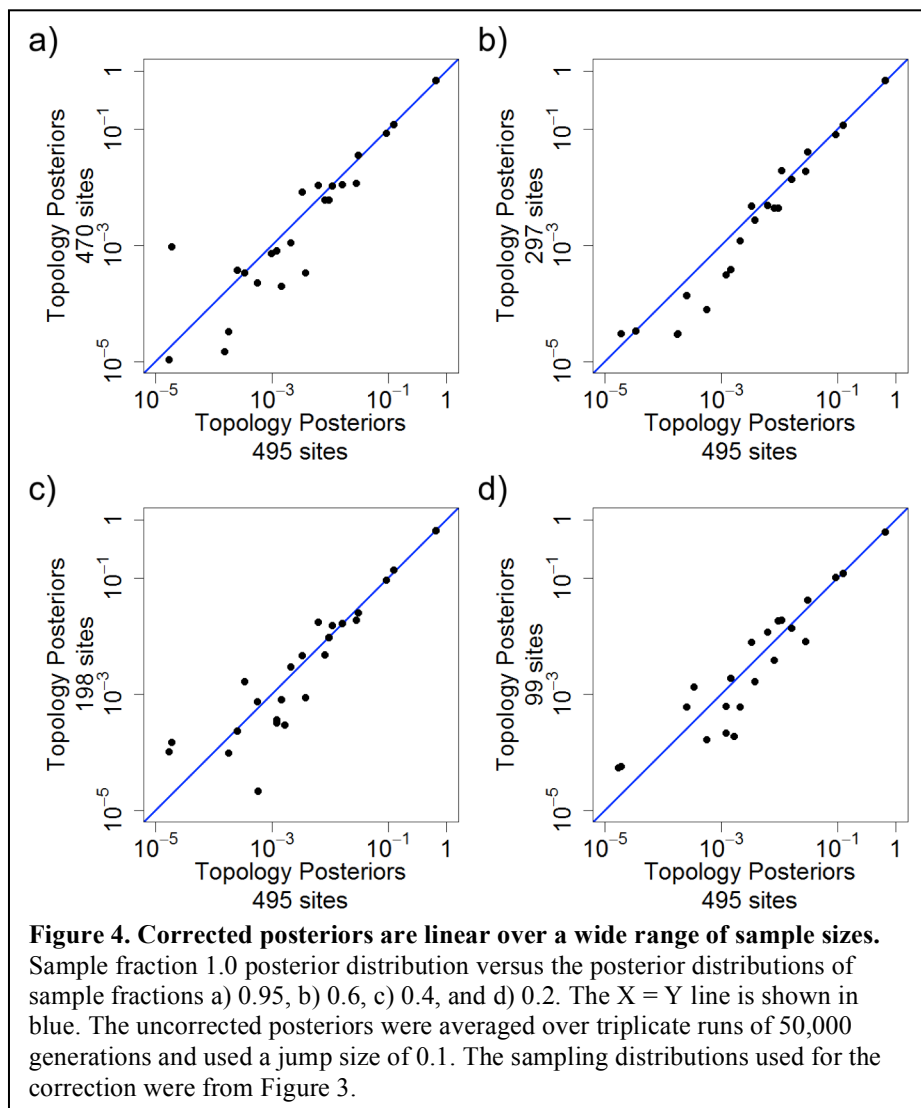
221 Predictably, comparisons
222 of replicate sampling
223 distribution runs indicated
224 an increasing variance in
225 estimated biases with
226 decreasing probabilities in
227 the runs (Fig. S2).

228 The posterior
229 correction (Equation 5)
230 appears to work well
231 across a broad range of
232 sample sizes (Fig. 4). The



MARKOV KATANA

233 corrected topology
234 posteriors for sample
235 fractions from 20%
236 to 95% were all
237 highly correlated
238 with the topology
239 posteriors for sample
240 size 100%. It should
241 be noted that the
242 uncorrected
243 posteriors are only
244 slightly less
245 correlated with each
246 other than are the
247 corrected posteriors
248 (Fig. S3), meaning



249 that the answer would have been similar without the correction. It is probably best to use the
250 correction anyway, because in more complicated situations it may make more of a difference,
251 and it is not too much trouble to obtain and is correct.

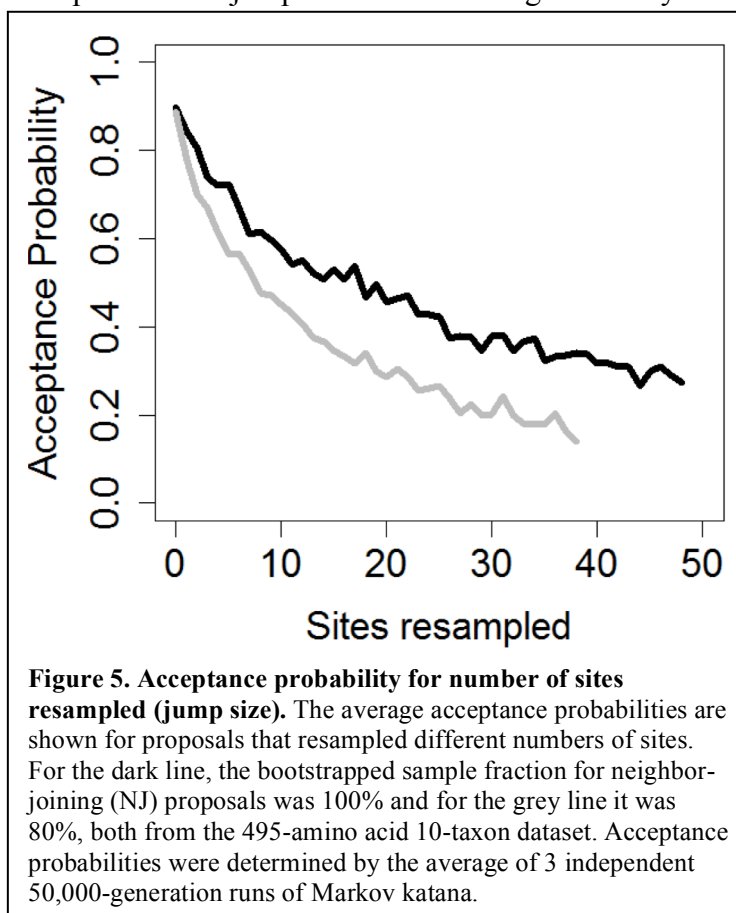
252 The corrected posterior estimates appear to be most noisy when the sampling distribution
253 estimate is small and therefore poorly estimated. This is not entirely surprising given that the
254 sampling distribution is in the denominator. Because the sampling distribution calculations are
255 computationally inexpensive (they do not require a likelihood calculation), it is possible to obtain

MARKOV KATANA

256 a couple orders of magnitude more data for them than for the uncorrected posterior estimates.
257 While estimating the sampling distribution more precisely is important for the correction, many
258 of the trees examined have topologies that are not found in the posterior. A potential means to
259 increase accuracy of relevant topologies in the sampling distribution is to limit the sampling prior
260 chains to those topologies seen in the uncorrected posterior.

261 *Effect of Sample Fraction and Jump Size on the Markov Chain*

262 Although the posterior estimates were comparable for all sample sizes, it is still
263 worthwhile to consider the effect of both sample size and jump size on the mixing efficiency of
264 the Markov chain. For a range of
265 conditions considered, the acceptance
266 probability for Markov chain proposal
267 varied from 10% to 90% (Fig. 5). We
268 chose 100% sample size bootstraps
269 for the NJ proposals along with a
270 jump size of 50 (10%) as standard
271 reference conditions, which had
272 acceptance probabilities of about
273 30%.



MARKOV KATANA

274 We also considered the effect of jump size on both the sampling distribution and the
275 uncorrected posterior Markov chain estimates. For the initial sampling distribution estimation
276 procedure, the most well-mixed chain is of course the one with independent bootstraps

277 ($j=100\%$), but the chain also

278 mixes well with lower jump

279 sizes. It is necessary to have

280 smaller jump sizes because a

281 high proportion (99.9%) of the

282 random samples are not in the

283 uncorrected posterior topology

284 set. For this analysis, the

285 optimal jump size was $j=0.85$.

286 This result did not differ much

287 for a range of sample sizes.

288 Although differing in detail,

289 the jump size analysis for the

290 uncorrected posterior had

291 similar results to the biased sampling prior analysis.

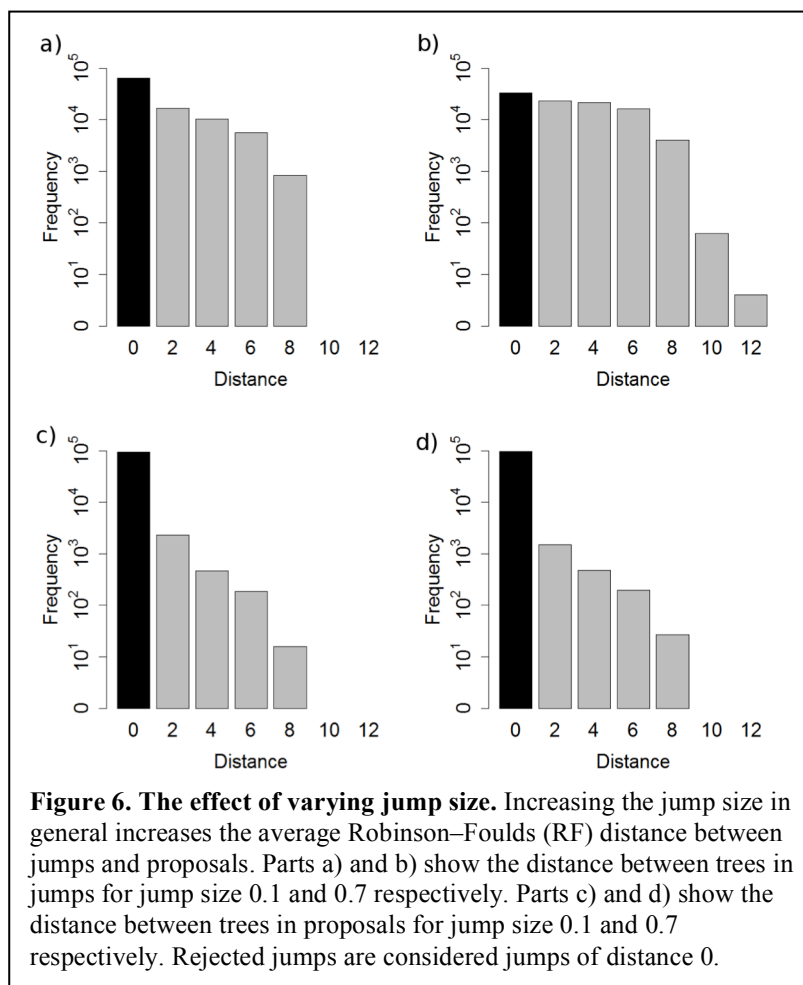
292 Acceptance probabilities varied widely depending on jump size. In general, 5-10 sites

293 appears to be a minimum, and 50 sites is probably a maximum. With smaller sample sizes (e.g.,

294 80% shown here), the jump size is a larger proportion of the sample and reduces the acceptance

295 probability more rapidly. Jump sizes bigger than 50 have somewhat greater probability of

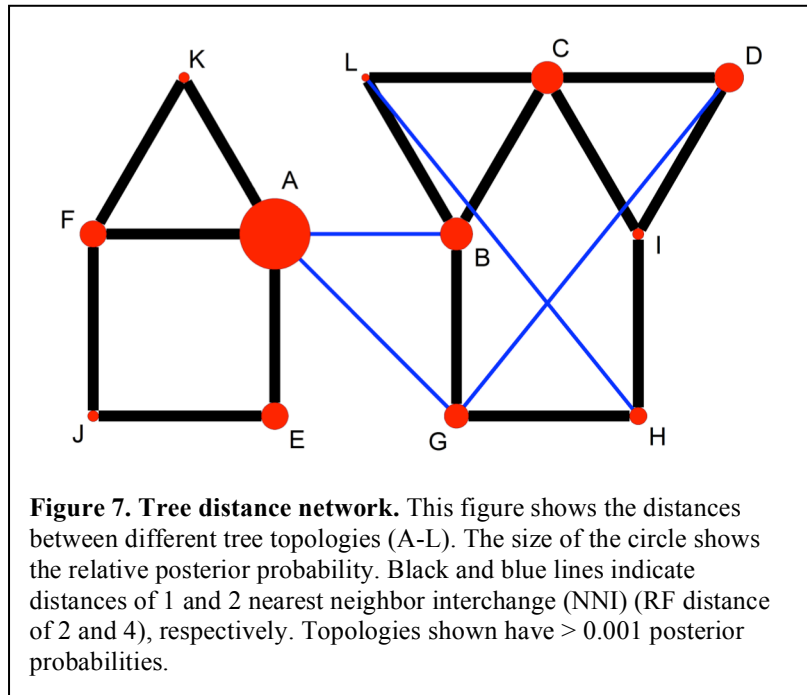
296 making large jumps in topology space (Fig. 6), at the cost of reduced probabilities of jumps to



297 the same topology due to lower acceptance probabilities.

298 *The Structure of Tree Space*

299 Figure 7 shows a
300 network representation of the
301 12 tree topologies that had a
302 posterior of 0.001 or higher
303 (see Table 2). The size of the
304 node represents the relative
305 posterior of the topology, and
306 the edges of the graph indicate
307 NNI distances of one or two
308 between the tree topologies.



309 The tree topology space of this test data was clearly divided into two clusters of trees shown by
310 the intragroup connections and the few intergroup connections. Given the connectivity of the
311 network, other tree topology sampling procedures may have difficulty jumping between groups.

312 DISCUSSION

313 We have demonstrated here that the Markov katana bootstrapping approach to
314 phylogenetic tree searching can be a highly effective means for finding Bayesian posterior
315 topologies and branches. It is able to take advantage of the speed of approximate distance-based
316 methods to propose new trees, but retains the reliability of Bayesian methods. Many previous
317 phylogenetic tree-search methods use the provided sequences for only the likelihood
318 calculations, but Markov katana introduces a new way to explore tree space informed by the

MARKOV KATANA

319 sequences. Including the sequence data in the tree search improves the fraction of high likelihood
320 trees proposed and allows efficient jump proposals between even distant topologies.

321 For the 10-taxon dataset, the NJ algorithm is extremely fast, and the overall speed of the
322 MK computation was limited by the likelihood calculations. As the number of taxa grows
323 beyond ~200, the NJ algorithm slows dramatically and dominates computation times (data not
324 shown). This could be alleviated using fast heuristic NJ algorithms or external programs such as
325 RapidNJ that are optimized for large alignments (Simonsen 2008). Our current implementation
326 calculates the distance contribution of each site only once and so is not hindered by the
327 complexity of the distance measure. We did not see a great difference in the proposal bias for the
328 two distance measures we compared, but further exploration of the performance of alternative
329 distances may in some cases be warranted.

330 We used PAML for the likelihood calculations, but any program that computes
331 likelihoods could potentially be used. The simplicity and adjustability of the approach means that
332 it could be easily incorporated into existing sequence analysis packages (e.g., MrBayes, PAUP*,
333 HyPhy, and PAML (Ronquist 2003; Swofford 2003; Pond 2005; Yang 2007)). We used a Perl
334 script to implement the MK algorithm and demonstrate the method as simply as possible, but we
335 expect that MK can be easily integrated directly into existing programs, which would then
336 undoubtedly be much faster. We did not see the benefit in constructing a new likelihood program
337 from scratch, although we believe the methodology would interact well with our existing
338 context-dependent Bayesian analysis program, *PLEX* (de Koning 2012).

339 ACKNOWLEDGEMENTS

340 Thanks to Seena D. Shah, who contributed to early versions of coding on *MarkovKatana*. We

MARKOV KATANA

341 acknowledge the support of the National Institutes of Health (NIH; GM083127 and
342 GM097251) to DDP.

343 LITERATURE CITED

- 344 Alfaro, Michael E., Stefan Zoller, and François Lutzoni. 2003. "Bayes or Bootstrap? A
345 Simulation Study Comparing the Performance of Bayesian Markov Chain Monte Carlo
346 Sampling and Bootstrapping in Assessing Phylogenetic Confidence." *Molecular Biology
347 and Evolution* 20 (2): 255–66. doi:10.1093/molbev/msg028.
- 348 Anisimova, Maria, and Olivier Gascuel. 2006. "Approximate Likelihood-Ratio Test for
349 Branches: A Fast, Accurate, and Powerful Alternative." *Systematic Biology* 55 (4): 539–
350 52. doi:10.1080/10635150600755453.
- 351 Bodlaender, Hans, Mike Fellows, and Tandy Warnow. 1992. "Two Strikes against Perfect
352 Phylogeny." *Automata, Languages and Programming*, 273–283.
- 353 Brocchieri, Luciano. 2001. "Phylogenetic Inferences from Molecular Sequences: Review and
354 Critique." *Theoretical Population Biology* 59 (1): 27–40. doi:10.1006/tpbi.2000.1485.
- 355 Castoe, T. A., Z. J. Jiang, W. Gu, Z. O. Wang, and D. D. Pollock. 2008. "Adaptive Evolution and
356 Functional Redesign of Core Metabolic Proteins in Snakes." *PLoS One* 3 (5): e2201.
357 doi:10.1371/journal.pone.0002201.
- 358 Castoe, Todd A., A. P. Jason de Koning, Hyun-Min Kim, Wanjun Gu, Brice P. Noonan, Gavin
359 Naylor, Zhi J. Jiang, Christopher L. Parkinson, and David D. Pollock. 2009. "Evidence
360 for an Ancient Adaptive Episode of Convergent Molecular Evolution." *Proceedings of
361 the National Academy of Sciences*, April, pnas.0900233106.
362 doi:10.1073/pnas.0900233106.
- 363 Craig, Roger A., and Li Liao. 2007. "Phylogenetic Tree Information Aids Supervised Learning
364 for Predicting Protein-Protein Interaction Based on Distance Matrices." *BMC
365 Bioinformatics* 8: 6. doi:10.1186/1471-2105-8-6.
- 366 Efron, Bradley, Elizabeth Halloran, and Susan Holmes. 1996. "Bootstrap Confidence Levels for
367 Phylogenetic Trees." *Proceedings of the National Academy of Sciences* 93 (23): 13429–
368 13429.
- 369 Felsenstein, Joseph. 1981. "Evolutionary Trees from DNA Sequences: A Maximum Likelihood
370 Approach." *Journal of Molecular Evolution* 17 (6): 368–76. doi:10.1007/BF01734359.
- 371 Felsenstein, Joseph. 1984. "Distance Methods for Inferring Phylogenies: A Justification."
372 *Evolution* 18 (1): 16–24. doi:10.2307/2408542.
- 373 Felsenstein, Joseph. 1985. "Confidence Limits on Phylogenies: An Approach Using the
374 Bootstrap." *Evolution* 39 (4): 783–91. doi:10.2307/2408678.
- 375 Felsenstein, Joseph. 2004. *Inferring Phylogenies*. Vol. 2. Sunderland, Massachusetts: Sinauer
376 Associates. <http://www.sinauer.com/media/wysiwyg/tocs/InferringPhylogenies.pdf>.
- 377 Fukushima, Kenji, Xiaodong Fang, David Alvarez-Ponce, Huimin Cai, Lorenzo Carretero-
378 Paulet, Cui Chen, Tien-Hao Chang, Kimberly M. Farr, Tomomichi Fujita, Yuji
379 Hiwatashi, and others. 2017. "Genome of the Pitcher Plant *Cephalotus* Reveals Genetic
380 Changes Associated with Carnivory." *Nature Ecology & Evolution* 1: 0059.
- 381 Geyer, Charles J. 1991. "Markov Chain Monte Carlo Maximum Likelihood." In . *Interface*

MARKOV KATANA

- 382 Foundation of North America. <http://conservancy.umn.edu/handle/11299/58440>.
- 383 Goldstein, Richard A., Stephen T. Pollard, Seena D. Shah, and David D. Pollock. 2015.
- 384 “Nonadaptive Amino Acid Convergence Rates Decrease over Time.” *Molecular Biology*
- 385 *and Evolution* 32 (6): 1373–81. doi:10.1093/molbev/msv041.
- 386 Hackett, Jeremiah D., Hwan Su Yoon, Shenglan Li, Adrian Reyes-Prieto, Susanne E. Rümmele,
- 387 and Debashish Bhattacharya. 2007. “Phylogenomic Analysis Supports the Monophyly of
- 388 Cryptophytes and Haptophytes and the Association of Rhizaria with Chromalveolates.”
- 389 *Molecular Biology and Evolution* 24 (8): 1702–1713.
- 390 Hastings, W. K. 1970. “Monte Carlo Sampling Methods Using Markov Chains and Their
- 391 Applications.” *Biometrika* 57 (1): 97–109.
- 392 Hendy, M. D., and David Penny. 1982. “Branch and Bound Algorithms to Determine Minimal
- 393 Evolutionary Trees.” *Mathematical Biosciences* 59 (2): 277–90. doi:10.1016/0025-
- 394 5564(82)90027-X.
- 395 Hershkovitz, Mark A., and Detlef D. Leipe. 1998. “Phylogenetic Analysis.” In *Bioinformatics*,
- 396 edited by Andreas D. Baxevanis and B. F. Francis Ouellette, 189–230. John Wiley &
- 397 Sons, Inc. doi:10.1002/9780470110607.ch9.
- 398 Huelsenbeck, John P. 1995. “Performance of Phylogenetic Methods in Simulation.” *Systematic*
- 399 *Biology* 44 (1): 17–48. doi:10.1093/sysbio/44.1.17.
- 400 Huelsenbeck, John P., and Keith A. Crandall. 1997. “Phylogeny Estimation and Hypothesis
- 401 Testing Using Maximum Likelihood.” *Annual Review of Ecology and Systematics* 28 (1):
- 402 437–66. doi:10.1146/annurev.ecolsys.28.1.437.
- 403 Huelsenbeck, John P., and Mark Kirkpatrick. 1996. “Do Phylogenetic Methods Produce Trees
- 404 with Biased Shapes?” *Evolution* 50 (4): 1418–24. doi:10.2307/2410879.
- 405 Huelsenbeck, John P., Fredrik Ronquist, and others. 2001. “MRBAYES: Bayesian Inference of
- 406 Phylogenetic Trees.” *Bioinformatics* 17 (8): 754–755.
- 407 Kirkpatrick, S., C. D. Gelatt, and M. P. Vecchi. 1983. “Optimization by Simulated Annealing.”
- 408 *Science* 220 (4598): 671–80.
- 409 Koning, A. P. Jason de, Wanjun Gu, Todd A. Castoe, and David D. Pollock. 2012.
- 410 “Phylogenetics, Likelihood, Evolution and Complexity.” *Bioinformatics* 28 (22): 2989–
- 411 90. doi:10.1093/bioinformatics/bts555.
- 412 Koning, AP Jason de, Wanjun Gu, and David D. Pollock. 2010. “Rapid Likelihood Analysis on
- 413 Large Phylogenies Using Partial Sampling of Substitution Histories.” *Molecular Biology*
- 414 *and Evolution* 27 (2): 249–265.
- 415 Kuhner, M. K., J. Yamato, and J. Felsenstein. 1995. “Estimating Effective Population Size and
- 416 Mutation Rate from Sequence Data Using Metropolis-Hastings Sampling.” *Genetics* 140
- 417 (4): 1421–30.
- 418 Matsuda, Hideo. 1995. “Construction of Phylogenetic Trees from Amino Acid Sequences Using
- 419 a Genetic Algorithm.” *Genome Informatics*, July, 19. doi:10.11234/gi1990.6.19.
- 420 Mossel, Elchanan, and Eric Vigoda. 2005. “Phylogenetic MCMC Algorithms Are Misleading on
- 421 Mixtures of Trees.” *Science* 309 (5744): 2207–9. doi:10.1126/science.1115493.
- 422 Nixon, K. 1999. “The Parsimony Ratchet, a New Method for Rapid Parsimony Analysis.”
- 423 *Cladistics* 15 (4): 407–14. doi:10.1006/clad.1999.0121.
- 424 Pickett, Kurt M., and Christopher P. Randle. 2005. “Strange Bayes Indeed: Uniform Topological
- 425 Priors Imply Non-Uniform Clade Priors.” *Molecular Phylogenetics and Evolution* 34 (1):
- 426 203–11. doi:10.1016/j.ympev.2004.09.001.
- 427 Pollock, David D. 1998. “Increased Accuracy in Analytical Molecular Distance Estimation.”

MARKOV KATANA

- 428 *Theoretical Population Biology* 54 (1): 78–90. doi:10.1006/tpbi.1998.1362.
- 429 Pollock, David D., and William J. Bruno. 2000. “Assessing an Unknown Evolutionary Process:
430 Effect of Increasing Site-Specific Knowledge Through Taxon Addition.” *Molecular*
431 *Biology and Evolution* 17 (12): 1854–58.
- 432 Pond, Sergei L. Kosakovsky, and Spencer V. Muse. 2005. “HyPhy: Hypothesis Testing Using
433 Phylogenies.” In *Statistical Methods in Molecular Evolution*, 125–181. Springer.
434 http://link.springer.com/content/pdf/10.1007/0-387-27733-1_6.pdf.
- 435 Reyes-Prieto, Adrian, and Debashish Bhattacharya. 2007. “Phylogeny of Nuclear-Encoded
436 Plastid-Targeted Proteins Supports an Early Divergence of Glaucophytes within Plantae.”
437 *Molecular Biology and Evolution* 24 (11): 2358–2361.
- 438 Ronquist, Fredrik, and John P. Huelsenbeck. 2003. “MrBayes 3: Bayesian Phylogenetic
439 Inference under Mixed Models.” *Bioinformatics* 19 (12): 1572–1574.
- 440 Saitou, N., and M. Nei. 1987. “The Neighbor-Joining Method: A New Method for
441 Reconstructing Phylogenetic Trees.” *Molecular Biology and Evolution* 4 (4): 406–25.
442 doi:10.1093/oxfordjournals.molbev.a040454.
- 443 Salemi, Marco, Philippe Lemey, and Anne-Mieke Vandamme. 2009. *The Phylogenetic*
444 *Handbook: A Practical Approach to Phylogenetic Analysis and Hypothesis Testing*.
445 Cambridge University Press.
- 446 Sanderson, Michael J., and Amy C. Driskell. 2003. “The Challenge of Constructing Large
447 Phylogenetic Trees.” *Trends in Plant Science* 8 (8): 374–79. doi:10.1016/S1360-
448 1385(03)00165-1.
- 449 Simonsen, Martin, Thomas Mailund, and Christian N. S. Pedersen. 2008. “Rapid Neighbour-
450 Joining.” In *Algorithms in Bioinformatics*, edited by Keith A. Crandall and Jens
451 Lagergren, 5251:113–22. Berlin, Heidelberg: Springer Berlin Heidelberg.
452 doi:10.1007/978-3-540-87361-7_10.
- 453 Sullivan, Jack, and Paul Joyce. 2005. “Model Selection in Phylogenetics.” *Annual Review of*
454 *Ecology, Evolution, and Systematics* 36 (1): 445–66.
455 doi:10.1146/annurev.ecolsys.36.102003.152633.
- 456 Swofford, D. L. 2003. *PAUP* Phylogenetic Analysis Using Parsimony (*and Other Methods)*
457 (version Version 4). Sunderland, Massachusetts: Sinauer Associates.
- 458 Takahashi, Kei, and Masatoshi Nei. 2000. “Efficiencies of Fast Algorithms of Phylogenetic
459 Inference Under the Criteria of Maximum Parsimony, Minimum Evolution, and
460 Maximum Likelihood When a Large Number of Sequences Are Used.” *Molecular*
461 *Biology and Evolution* 17 (8): 1251–58. doi:10.1093/oxfordjournals.molbev.a026408.
- 462 Vonk, Freek J., Nicholas R. Casewell, Christiaan V. Henkel, Alysha M. Heimberg, Hans J.
463 Jansen, Ryan J. R. McCleary, Harald M. E. Kerckamp, Rutger A. Vos, Isabel Guerreiro,
464 Juan J. Calvete, Wolfgang Wüster, Anthony E. Woods, Jessica M. Logan, Robert A.
465 Harrison, Todd A. Castoe, A. P. Jason de Koning, David D. Pollock, Mark Yandell,
466 Diego Calderon, Camila Renjifo, Rachel B. Currier, David Salgado, Davinia Pla, Libia
467 Sanz, Asad S. Hyder, José M. C. Ribeiro, Jan W. Arntzen, Guido E. E. J. M. van den
468 Thillart, Marten Boetzer, Walter Pirovano, Ron P. Dirks, Herman P. Spaink, Denis
469 Duboule, Edwina McGlinn, R. Manjunatha Kini, and Michael K. Richardson. 2013. “The
470 King Cobra Genome Reveals Dynamic Gene Evolution and Adaptation in the Snake
471 Venom System.” *Proceedings of the National Academy of Sciences* 110 (51): 20651–56.
472 doi:10.1073/pnas.1314702110.
- 473 Vos, R. A. 2003. “Accelerated Likelihood Surface Exploration: The Likelihood Ratchet.”

MARKOV KATANA

- 474 *Systematic Biology* 52 (3): 368–73. doi:10.1080/10635150390196993.
- 475 Wang, Zhengyuan O., and David D. Pollock. 2005. “Context Dependence and Coevolution
- 476 among Amino Acid Residues in Proteins.” *Methods in Enzymology* 395: 779–790.
- 477 Whelan, Simon. 2007. “New Approaches to Phylogenetic Tree Search and Their Application to
- 478 Large Numbers of Protein Alignments.” *Systematic Biology* 56 (5): 727–40.
- 479 doi:10.1080/10635150701611134.
- 480 Whelan, Simon, Pietro Liò, and Nick Goldman. 2001. “Molecular Phylogenetics: State-of-the-
- 481 Art Methods for Looking into the Past.” *Trends in Genetics* 17 (5): 262–72.
- 482 doi:10.1016/S0168-9525(01)02272-7.
- 483 Xia, Xuhua. 2006. “Topological Bias in Distance-Based Phylogenetic Methods: Problems with
- 484 over-and Underestimated Genetic Distances.” *Evolutionary Bioinformatics* 2.
- 485 [http://search.proquest.com/openview/186314449a175c3222c8c80061e94530/1?pq-](http://search.proquest.com/openview/186314449a175c3222c8c80061e94530/1?pq-origsite=gscholar&cbl=1026404)
- 486 [origsite=gscholar&cbl=1026404.](http://search.proquest.com/openview/186314449a175c3222c8c80061e94530/1?pq-origsite=gscholar&cbl=1026404)
- 487 Yang, Ziheng. 2007. “PAML 4: Phylogenetic Analysis by Maximum Likelihood.” *Mol Biol Evol*
- 488 24 (8): 1586–91. doi:10.1093/molbev/msm088.
- 489 Yang, Ziheng, R. Nielsen, and M. Hasegawa. 1998. “Models of Amino Acid Substitution and
- 490 Applications to Mitochondrial Protein Evolution.” *Mol Biol Evol* 15 (12): 1600–1611.
- 491 Yang, Ziheng, and Bruce Rannala. 2005. “Branch-Length Prior Influences Bayesian Posterior
- 492 Probability of Phylogeny.” *Systematic Biology* 54 (3): 455–70.
- 493 doi:10.1080/10635150590945313.
- 494 Zharkikh, Andrey, and Wen-Hsiung Li. 1995. “Estimation of Confidence in Phylogeny: The
- 495 Complete-and-Partial Bootstrap Technique.” *Molecular Phylogenetics and Evolution* 4
- 496 (1): 44–63. doi:10.1006/mpev.1995.1005.

497

498

499 FIGURE CAPTIONS

500 **Figure 1. Flow of the Markov katana procedure.**

501
502 **Figure 2. An example phylogenetic tree for the 10-taxon dataset.** A tree for all protein coding regions in the
503 mitochondrial genome is shown. Posterior probabilities are shown for the 495 amino acid Cytochrome C Oxidase
504 subunit 1 (COI) alignment used in initial analyses. The posteriors were calculated using Mr. Bayes. Given the
505 limited data used for the purpose of testing our method, this tree should not be interpreted as a true or species tree.

506
507 **Figure 3. Relative sampling distributions estimations comparing different sample sizes.** Sample size 1.0
508 sampling distribution versus the sampling distributions of sample size a) 0.95, b) 0.6, c) 0.4, d) 0.2. The blue line
509 indicates where x and y values are equal. The sampling distributions were averaged over triplicate runs of 1,000,000
510 generations and used a jump size of 0.1.

511
512 **Figure 4. Corrected posteriors are linear over a wide range of sample sizes.** Sample fraction 1.0 posterior
513 distribution versus the posterior distributions of sample fractions a) 0.95, b) 0.6, c) 0.4, and d) 0.2. The X = Y line is
514 shown in blue. The uncorrected posteriors were averaged over triplicate runs of 50,000 generations and used a jump
515 size of 0.1. The sampling distributions used for the correction were from Figure 3.

516
517 **Figure 5. Acceptance probability for number of sites resampled (jump size).** The average acceptance
518 probabilities are shown for proposals that resampled different numbers of sites. For the dark line, the bootstrapped
519 sample fraction for neighbor-joining (NJ) proposals was 100% and for the grey line it was 80%, both from the 495-
520 amino acid 10-taxon dataset. Acceptance probabilities were determined by the average of 3 independent 50,000
521 generation runs of Markov katana.

522
523 **Figure 6. The effect of varying jump size.** Increasing the jump size in general increases the average Robinson-
524 Foulds (RF) distance between jumps and proposals. Parts a) and b) show the distance between trees in jumps for
525 jump size 0.1 and 0.7 respectively. Parts c) and d) show the distance between trees in proposals for jump size 0.1
526 and 0.7 respectively. Rejected jumps are considered jumps of distance 0.

527
528 **Figure 7. Tree distance network.** This figure shows the distances between different tree topologies (A-L). The size
529 of the circle shows the relative posterior probability. Black and blue lines indicate distances of 1 and 2 nearest
530 neighbor interchange (NNI) (RF distance of 2 and 4), respectively. Topologies that are > 0.001 posterior probability.

531

532 TABLES

533 **Table 1: Species in the Cytochrome C Oxidase Subunit 1 Alignment**

<i>Chrysochloris asiatica</i>
<i>Monodelphis domestica</i>
<i>Chamaeleo calcaricarenis</i>
<i>Trachypithecus obscurus</i>
<i>Gulo gulo</i>
<i>Ensatina eschscholtzii</i>
<i>Cervus nippon centralis</i>
<i>Presbytis melalophos</i>

Anomalurus sp. GP-2005

Agapornis roseicollis

534

535 **Table 2: Tree Posteriors**

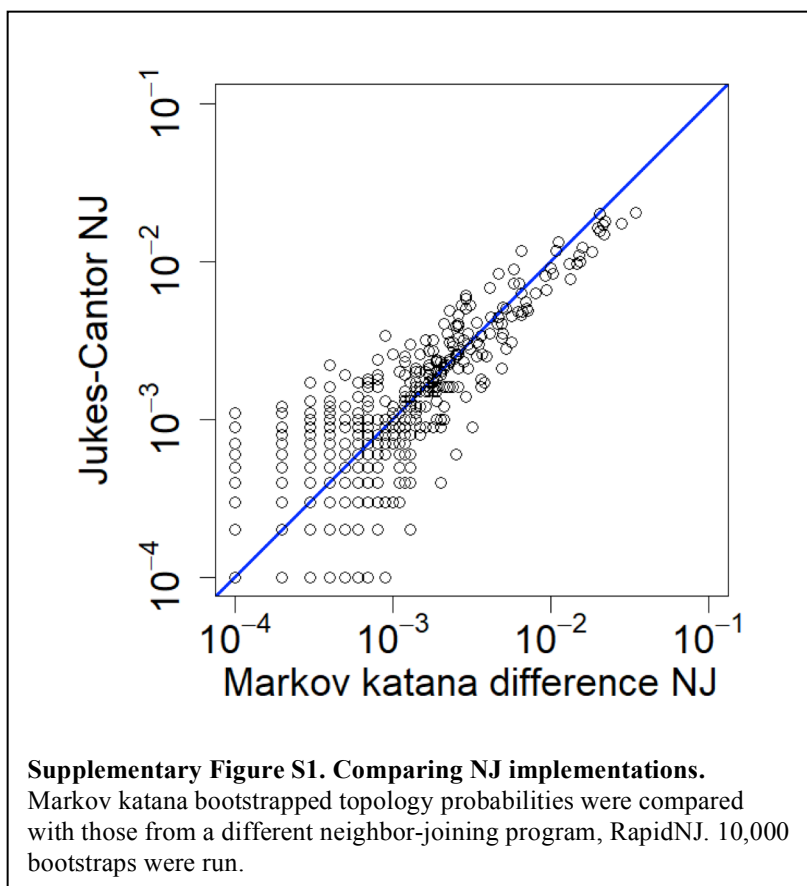
Letter	Id	Posterior
A	16405	0.692
B	78835	0.067
C	36915	0.066
D	92575	0.048
E	26545	0.038
F	57985	0.037
G	80055	0.029
H	39655	0.010
I	2955	0.003
J	82665	0.003
K	8085	0.002
L	88595	0.001

536 **Table 2.** Posterior probability for topologies with substantial representation in the
537 uncorrected posterior for the 10-taxon dataset. The topologies are labeled for reference in Figure
538 7.

539

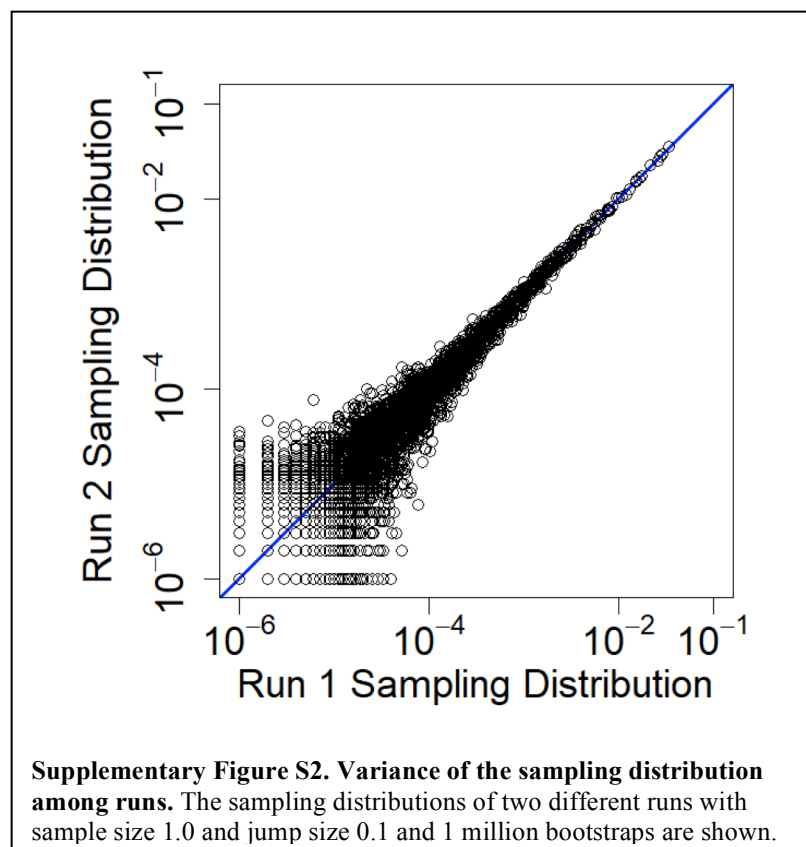
540

MARKOV KATANA



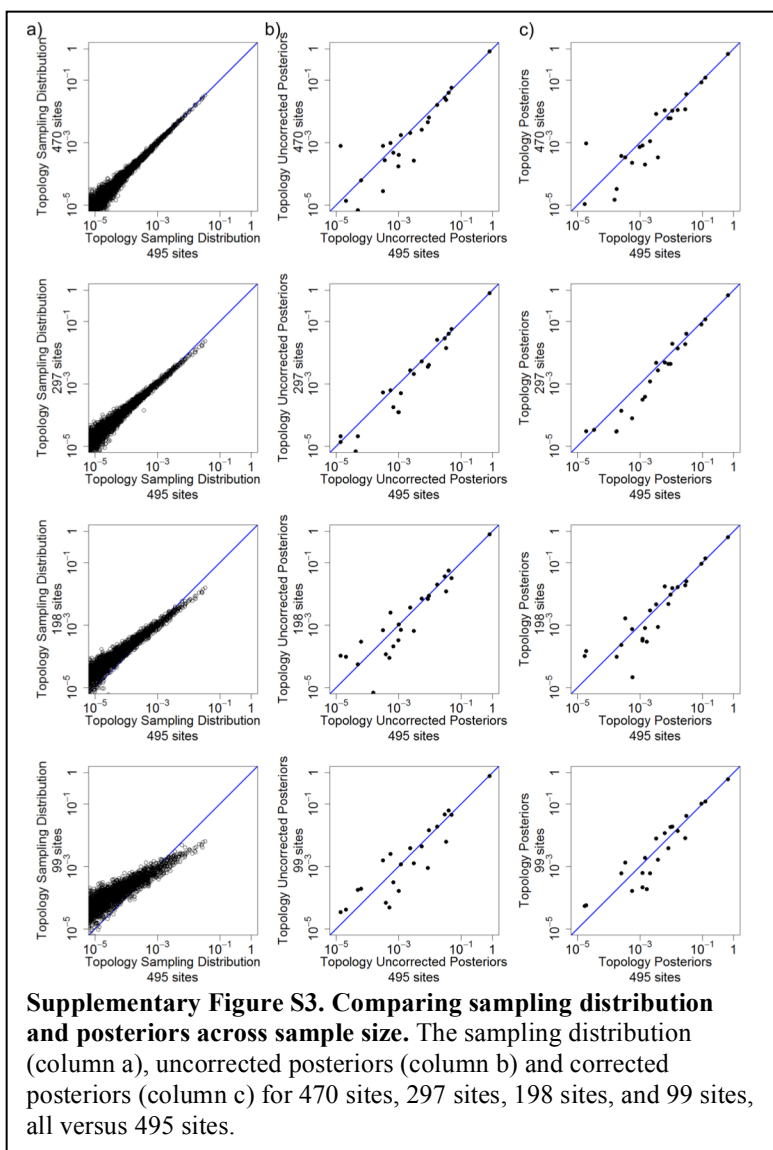
541

MARKOV KATANA



542

543



544

545